

AN UPPER BOUND ON THE CONVERGENCE RATE OF A SECOND FUNCTIONAL IN OPTIMAL SEQUENCE ALIGNMENT

RAPHAEL HAUSER, HEINRICH MATZINGER, AND IONEL POPESCU

ABSTRACT. Consider finite sequences $X_{[1,n]} = X_1 \dots X_n$ and $Y_{[1,n]} = Y_1 \dots Y_n$ of length n , consisting of i.i.d. samples of random letters from a finite alphabet, and let S and T be chosen i.i.d. randomly from the unit ball in the space of symmetric scoring functions over this alphabet augmented by a gap symbol. We prove a probabilistic upper bound of linear order in $n^{0.75}$ for the deviation of the score relative to T of optimal alignments with gaps of $X_{[1,n]}$ and $Y_{[1,n]}$ relative to S . It remains an open problem to prove a lower bound. Our result contributes to the understanding of the microstructure of optimal alignments relative to one given scoring function, extending a theory begun in [4].

1. INTRODUCTION AND MAIN RESULTS

The subject of this paper is concerned with the asymptotics of optimal sequence alignments for random sequences whose lengths tend to infinity. An important problem that occurs both in bioinformatics and in natural language processing is to decide on the homology of two (or more) finite sequences consisting of symbols from a fixed finite alphabet. A highly successful approach is to fix a *scoring function* and maximise the total score over the set of all alignments with gaps of the two sequences (for a precise definition, see the text below). Despite the combinatorially many alignments to be considered, the total score can be maximised in polynomial time by use of a dynamic

1991 *Mathematics Subject Classification.* Primary 60F10; Secondary 92D20, 60K35.

Key words and phrases. Sequence alignment, convex geometry, large deviations, percolation theory.

Raphael Hauser was supported by the Engineering and Physical Sciences Research Council [grant number EP/H02686X/1].

Ionel Popescu was partially supported by a grant of the Romanian National Authority for Scientific Research, CNCS - UEFISCDI, project number PN-II-RU-TE-2011-3-0259 and Marie Curie Action Grant PIRG.GA.2009.249200.

programming recursion [5]. Using this approach, two sequences can be considered as homologous if the total score of their optimal alignment relative to a salient scoring function significantly exceeds the typical total score of an optimal alignment of two random sequences of the same length. Rigorous statistical tests on this basis require an understanding of relevant null models, thus giving the initial motivation for the theoretical study of optimal sequence alignments of random sequences and their total scores [6].

The purpose of this paper is to contribute to this theory by studying the following question: given two *symmetric* scoring functions S and T , and given two i.i.d. random sequences of length n , does the rescaled total score (the score divided by n) relative to T of an optimal alignment of the two sequences relative to S converge as n tends to infinity, and if the answer to this question is ‘yes’, can we bound the convergence rate? We will answer both questions in the affirmative. Before we go into the technical details of our analysis, we introduce the necessary notation and background and give further details on the main contributions of this paper in relation to the existing literature.

1.1. Alignments with Gaps. Let $n \in \mathbb{N}$ and write $[1, n] := \{1, \dots, n\}$. Consider two sequences of length n , $x_{[1, n]} := (x_i)_{i \in [1, n]}$ and $y_{[1, n]} := (y_j)_{j \in [1, n]}$ consisting of letters from a finite alphabet \mathcal{A} . Let us augment this alphabet by a symbol G for a *gap* and write $\mathcal{A}^* = \mathcal{A} \cup \{G\}$. We define an *alignment* (with gaps) of $x_{[1, n]}$ and $y_{[1, n]}$ as a pair of increasing subsequences $(i_\ell)_{\ell \in [1, k]}$ and $(j_\ell)_{\ell \in [1, k]}$ of $[1, n]$. For $\ell \in [1, k]$, each letter x_{i_ℓ} of the first sequence is then interpreted as aligned with the letter y_{j_ℓ} from the second sequence, while all remaining letters of either sequence are thought of as aligned with gaps.

For example the pair of increasing subsequences $(\{1, 5, 6, 8\}, \{2, 4, 5, 6\})$ of $[1, 8]$ correspond to the alignment

$$\begin{array}{cccccccccccc} G & x_1 & x_2 & x_3 & x_4 & G & x_5 & x_6 & x_7 & x_8 & G & G \\ y_1 & y_2 & G & G & G & y_3 & y_4 & y_5 & G & y_6 & y_7 & y_8 \end{array}$$

Note that the same subsequences also correspond to the alignment

$$\begin{array}{cccccccccccc} G & x_1 & x_2 & G & x_3 & x_4 & x_5 & x_6 & x_7 & x_8 & G & G \\ y_1 & y_2 & G & y_3 & G & G & y_4 & y_5 & G & y_6 & y_7 & y_8 \end{array}$$

and other arrangements obtained by permuting the order of consecutive letters aligned with gaps, so that the pair $(\{1, 5, 6, 8\}, \{2, 4, 5, 6\})$ represent in fact an equivalence class of alignments. By slight abuse of language, we will speak about an *alignment* when in fact referring to an entire equivalence class. In order to refer to the set of alignments of two sequences of length n , we introduce the following notation,

$$\Lambda_{n,k} := \left\{ ((i_\ell)_{\ell \in [1,k]}, (j_\ell)_{\ell \in [1,k]}) : 1 \leq i_1 < \dots < i_k \leq n, 1 \leq j_1 < \dots < j_k \leq n \right\}, \quad (k \in [0, n]),$$

$$\Lambda_n := \bigcup_{k=0}^n \Lambda_{n,k}.$$

1.2. Scoring Functions and Optimal Alignments. A function $R : \mathcal{A}^* \times \mathcal{A}^* \rightarrow \mathbb{R}$ will be called a *symmetric scoring function* if $R(\alpha, \beta) = R(\beta, \alpha)$ for all $\alpha, \beta \in \mathcal{A}^*$, and $R(G, G) = 0$. Given a symmetric scoring function R and two finite sequences $x_{[1,n]}$ and $y_{[1,n]}$ consisting of letters from the alphabet \mathcal{A} , we define the *total score* of $x_{[1,n]}$ and $y_{[1,n]}$ under an alignment $\nu = ((i_\ell), (j_\ell)) \in \Lambda_{n,k}$ as the sum of the scores of individually aligned letter pairs,

$$R_\nu(x_{[1,n]}, y_{[1,n]}) := \sum_{\ell=1}^k R(x_{i_\ell}, y_{j_\ell}) + \sum_{i \in [1,n] \setminus \{i_\ell : \ell \in [1,k]\}} R(x_i, G) + \sum_{j \in [1,n] \setminus \{j_\ell : \ell \in [1,k]\}} R(G, y_j).$$

Note that since our definition of alignments with gaps disallows the situation where a gap is aligned with a gap, the value of $R(G, G)$ should be inconsequential. Our rationale for requiring $R(G, G) = 0$ is to simplify some of our formulas, notably the norms defined in Section 1.5.

The *optimal alignment score* of $x_{[1,n]}$ and $y_{[1,n]}$ relative to R is defined by

$$R^*(x_{[1,n]}, y_{[1,n]}) := \max_{\nu \in \Lambda_n} R_\nu(x_{[1,n]}, y_{[1,n]}),$$

while the set of *optimal alignments* of $x_{[1,n]}$ and $y_{[1,n]}$ relative to R is the set of alignments

$$\nu_R^*(x_{[1,n]}, y_{[1,n]}) := \left\{ \nu \in \Lambda_n : R_\nu(x_{[1,n]}, y_{[1,n]}) = R^*(x_{[1,n]}, y_{[1,n]}) \right\}$$

on which the maximum is achieved. Note that in general, ν^* is not a singleton.

1.3. Random Sequences. Let us now consider two sequences $(X_i)_{i \in \mathbb{N}} : \Omega \rightarrow \mathcal{A}^{\mathbb{N}}$ and $(Y_j)_{j \in \mathbb{N}} : \Omega \rightarrow \mathcal{A}^{\mathbb{N}}$, defined on some appropriate probability space (Ω, \mathcal{F}, P) so as to consist of i.i.d. random letters X_i (respectively Y_i) drawn from a fixed probability distribution over a finite alphabet \mathcal{A} . Let us again augment this alphabet by a symbol G for a *gap* and write $\mathcal{A}^* = \mathcal{A} \cup \{G\}$. We write $X_{[1,n]} = (X_i)_{i=1}^n$ for the finite sequence consisting of the first n terms of $(X_i)_{i \in \mathbb{N}}$ and use a similar notation for the second sequence.

Let a symmetric scoring function R be given on $\mathcal{A}^* \times \mathcal{A}^*$. The following is then a well defined random variable for any $n \in \mathbb{N}$

$$L_{n,R} : \Omega \rightarrow \mathbb{R},$$

$$\omega \mapsto R^*(X_{[1,n]}(\omega), Y_{[1,n]}(\omega)),$$

and we write

$$\nu_{n,R}^* : \Omega \rightarrow \mathcal{P}(\Lambda_n),$$

$$\omega \mapsto \nu_R^*(X_{[1,n]}(\omega), Y_{[1,n]}(\omega))$$

for the random set of optimal alignments of $X_{[1,n]}$ and $Y_{[1,n]}$ relative to R .

It was shown in [2] that

$$(1.1) \quad \frac{L_{n,R}}{n} \xrightarrow{n \rightarrow \infty} \lambda_R \quad \text{almost surely,}$$

where λ_R is some deterministic constant that depends only on R . In Lemma 2.1 we give a proof that also establishes a quantitative convergence bound.

1.4. The Problem Setting of this Paper. Let us now consider two different symmetric scoring functions S and T and investigate the total score relative to T of an optimal alignment relative to S . Using the random sequences introduced above, we define the following random subsets of \mathbb{R}^2 ,

$$\text{SCORES}_{S,T}^n := \left\{ \left(\frac{S_\nu(X_{[1,n]}, Y_{[1,n]})}{n}, \frac{T_\nu(X_{[1,n]}, Y_{[1,n]})}{n} \right) : \nu \in \Lambda_n \right\}$$

$$\text{SET}_{S,T}^n := \text{cl}(\text{conv}(\text{SCORES}_{S,T}^n)),$$

where $\text{cl}(\cdot)$ denotes the topological closure in the canonical topology of \mathbb{R}^2 and $\text{conv}(\cdot)$ denotes the convex hull.

Next, consider a symmetric scoring function $R = aS + bT$ given as a linear combination of S and T . It follows from our definition of $\text{SET}_{S,T}^n$ that

$$(1.2) \quad \frac{L_{n,R}}{n} = \max_{(x,y) \in \text{SET}_{S,T}^n} f_{(a,b)}(x,y),$$

where $f_{(a,b)} : (x,y) \mapsto ax + by$ is the linear form on \mathbb{R}^2 defined by the weights a, b . Combining Equations (1.1) and (1.2), it follows that

$$(1.3) \quad \max_{(x,y) \in \text{SET}_{S,T}^n} f_{(a,b)}(x,y) \xrightarrow{n \rightarrow \infty} \lambda_{aS+bT}, \quad \text{a.s..}$$

We observe that, if a sequence of random compact convex sets $A_1, A_2, \dots \subset \mathbb{R}^2$ has the property that for any linear functional $f \in (\mathbb{R}^2)^*$,

$$\max_{(x,y) \in A_n} f(x,y) \xrightarrow{n \rightarrow \infty} \xi_f, \quad \text{a.s.,}$$

where $\xi_f \in \mathbb{R}$ is a deterministic constant that depends only on f , then the sequence $(A_n)_{n \in \mathbb{N}}$ converges in Hausdorff distance to a convex compact set A . We will prove this claim in Lemma 2.4. For compact sets $A, B \subset \mathbb{R}^2$, the Hausdorff distance is defined as

$$(1.4) \quad d_H(A, B) = \max\{\sup_{x \in A} \inf_{y \in B} d(x,y), \sup_{y \in B} \inf_{x \in A} d(x,y)\},$$

where $d(x,y) = \|x - y\|_2$ denotes the Euclidean distance.

Equation (1.3) and the fact that any $f \in (\mathbb{R}^2)^*$ is of the form $f_{(a,b)}$ for some $(a,b) \in \mathbb{R}^2$ show that the above made observation is applicable to the sequence of sets $(\text{SET}_{S,T}^n)_{n \in \mathbb{N}}$. There exists therefore a deterministic convex compact set $\text{SET}_{S,T}$ for which

$$(1.5) \quad d_H(\text{SET}_{S,T}^n, \text{SET}_{S,T}) \xrightarrow{n \rightarrow \infty} 0, \quad \text{a.s.}$$

One of our goals is to refine this analysis and quantify an upper-bound on the rate of convergence. An upper bound on the convergence was given in [4] for scoring functions that are not necessarily symmetric. In this paper we give a much simpler proof that is made possible by exploiting the symmetry of scoring functions. Since most scoring functions used in applications are symmetric, the simplification is of interest.

Another goal is to study how much the total score relative to T varies when two random strings are aligned optimally relative to S . Note that we have

$$L_{n,S} = \max_{(x,y) \in \text{SET}_{S,T}^n} x.$$

In general, we should not expect that $\nu_{n,S}^*$ to be a singleton. In other words, there may exist multiple optimal alignments of $X_{[1,n]}$ and $Y_{[1,n]}$ relative to S . Therefore, we need to consider the following quantities,

$$(1.6) \quad \max_{\pi \in \nu_{n,S}^*} \frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} = \max \left\{ y : (x, y) \in \text{SET}_{S,T}^n, x = \frac{L_{n,S}}{n} \right\},$$

$$(1.7) \quad \min_{\pi \in \nu_{n,S}^*} \frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} = \min \left\{ y : (x, y) \in \text{SET}_{S,T}^n, x = \frac{L_{n,S}}{n} \right\},$$

$$(1.8)$$

Lemma 2.5 will establish that if $\max_{(x,y) \in \text{SET}_{S,T}} x$ has a unique maximiser (x_0, y_0) , then the upper and lower bounds (1.6), (1.7) both converge to y_0 almost surely.

1.5. Statement of the Main Results. To state the main results of this paper, we introduce the following norms on the set of symmetric scoring functions $R : \mathcal{A}^* \times \mathcal{A}^* \rightarrow \mathbb{R}$,

$$(1.9) \quad |R| := \max_{a,b,c \in \mathcal{A}^*} |R(a,b) - R(a,c)|, \quad (\text{the change norm}),$$

$$(1.10) \quad |R|_2 := \sqrt{\sum_{a,b \in \mathcal{A}^*} R^2(a,b)}, \quad (\text{the Frobenius norm}).$$

The change norm plays the following important role: given two finite sequences and a fixed alignment with gaps, changing a single letter of one of the two sequences into an arbitrary other letter from the alphabet \mathcal{A} changes the total score of the alignment by at most $|R|$.

Theorem 1.1. *Let S and T be two symmetric scoring functions on $\mathcal{A}^* \times \mathcal{A}^*$ such that the optimisation problem $\max_{(x,y) \in \text{SET}_{S,T}^n} x$ has a unique maximiser (x_0, y_0) and the boundary of $\text{SET}_{S,T}$ has curvature at least $k > 0$ at this point, then the following bound applies for large enough n , where e is the Euler constant,*

$$\mathbb{P} \left[\left| \frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - y_0 \right| \leq \frac{5|T| + 2\sqrt{30|S|}}{k} \left(\frac{\ln(n e)}{n} \right)^{1/4}, \quad \forall \pi \in \nu_{n,S}^* \right] \geq 1 - 3n^{-\ln n}.$$

In particular if both S and T have change norm less than 1, the statement of Theorem 1.1 simplifies to

$$P \left[\left| \frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - y_0 \right| \leq \frac{11}{k} \left(\frac{\ln(ne)}{n} \right)^{1/4}, \quad \forall \pi \in \nu_{n,S}^* \right] \geq 1 - 3n^{-\ln n}, \quad \forall n \gg 1.$$

The curvature condition at the point (x_0, y_0) means that one can parametrize the boundary $\partial \text{SET}_{S,T}$ of the set $\text{SET}_{S,T}$ by a curve $c(t)$ for t in a neighbourhood of 0, with $c(0) = (x_0, y_0)$ and $\|\dot{c}\|_2 = 1$ for all t , where \dot{c} denotes the derivative with respect to t , the curvature

$$\kappa(\partial \text{SET}_{S,T}, (x_0, y_0)) := \|\ddot{c}(0)\|_2$$

then being defined as the standard curvature of this curve at $t = 0$. By convention, we define the curvature at vertices of $\partial \text{SET}_{S,T}$ (points on the boundary where $\text{SET}_{S,T}$ has a normal cone with nonempty interior) to be $+\infty$. We postpone the proof of Theorem 1.1 until Section 3.

While Theorem 1.1 establishes that if the boundary of $\text{SET}_{S,T}$ has positive curvature at (x_0, y_0) , then the T -score on an S -optimal alignment has a fluctuation of order at most $O([\ln(n)/n]^{0.25})$, the conditions of this result are difficult to verify in practice. However, as the following result shows, they apply generically:

Theorem 1.2. *Let S and T be chosen i.i.d. uniformly at random from the Frobenius-unit sphere in the space of symmetric scoring functions. Then the following hold true,*

- (1) $\max_{(x,y) \in \text{SET}_{S,T}^n} x$ has a unique maximiser (x_0, y_0) almost surely,
- (2) for any real number $k > 0$,

$$P[\kappa(\partial \text{SET}_{S,T}, (x_0, y_0)) < k] \leq \frac{4k}{\pi},$$

where $\kappa(\partial \text{SET}_{S,T}, (x_0, y_0))$ is the curvature at (x_0, y_0) of the boundary of $\text{SET}_{S,T}$.

Combining Theorems 1.1 and 1.2, we arrive at the following conclusion:

Corollary 1.1. *If the symmetric scoring functions S and T are chosen as in Theorem 1.2, then almost surely there exists $k > 0$ such that*

$$P \left[\left| \frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - y_0 \right| \leq \frac{11}{\max(k, 1)} \left(\frac{\ln(ne)}{n} \right)^{1/4}, \quad \forall \pi \in \nu_{n,S}^* \right] \geq 1 - 3n^{-\ln n}, \quad \forall n \gg 1.$$

2. PRELIMINARY RESULTS AND THEIR PROOFS

In this section we derive the main estimates on which the proofs of our main theorems rely. We begin by giving the classical Azuma-Hoeffding – McDiarmid Inequality.

Theorem 2.1. *Let W_1, \dots, W_n i.i.d. random variables that take values in some set D , let $a > 0$ be a constant and $f : D^n \rightarrow \mathbb{R}$ a n -variate real function with the property that for any $i \in [1, n]$, $w \in D^n$ and $z \in D$,*

$$|f(w_1, w_2, \dots, w_n) - f(w_1, w_2, \dots, w_{i-1}, z, w_{i+1}, \dots, w_n)| \leq a.$$

Then, for any $\epsilon > 0$, the following inequalities hold true,

$$\mathbb{P}[|f(W_1, W_2, \dots, W_n) - \mathbb{E}[f(W_1, \dots, W_n)]| \geq \epsilon n] \leq 2 \exp(-\epsilon^2 n / (2a^2)),$$

$$\mathbb{P}[f(W_1, W_2, \dots, W_n) - \mathbb{E}[f(W_1, \dots, W_n)] \geq \epsilon n] \leq \exp(-\epsilon^2 n / (2a^2)).$$

For a proof, see e.g. [1].

Lemma 2.1. *For any symmetric scoring function $R : \mathcal{A}^* \times \mathcal{A}^* \rightarrow \mathbb{R}$ there exists a deterministic constant λ_R such that*

$$\frac{L_{n,R}}{n} \xrightarrow{n \rightarrow \infty} \lambda_R, \quad a.s.$$

Proof. It is trivial to see that the function $n \mapsto \mathbb{E}[L_{n,R}]$ is superadditive. Therefore and since the scoring function is bounded, we have

$$(2.1) \quad \mathbb{E}[L_{n,R}]/n \xrightarrow{n \rightarrow \infty} \lambda_R := \sup_{n \geq 1} \mathbb{E}[L_{n,R}]/n,$$

where $\sup_{n \geq 1} \mathbb{E}[L_{n,R}]/n$ is well defined. For any $\epsilon > 0$, let $D_{n,R}(\epsilon)$ denote the event

$$D_{n,R}(\epsilon) = \{|L_{n,R} - \mathbb{E}[L_{n,R}]| \geq \epsilon \ln(n) \sqrt{n}\}.$$

Applying Theorem 2.1 with $a = |R|$, we obtain

$$(2.2) \quad \mathbb{P}[D_{n,R}(\epsilon) \leq 2 \exp(-\epsilon^2 (\ln n)^2 / 2|R|^2)] = 2n^{-\frac{\epsilon^2 \ln n}{2|R|^2}}.$$

By virtue of Borel-Cantelli, the finite summability of (2.2) implies that almost surely at most a finite number of the events $D_{n,R}(\epsilon)$ will hold. Combined with (2.1), and using the fact that $\epsilon > 0$ was arbitrary, this implies the claim. \square

The next result gives the rate of convergence for of $E[L_n(R)]/n$ toward λ_R . A bound for non-symmetric scoring functions was given in [4]. Here we exploit the symmetry of R to give a tighter bound that we will use to prove our main theorems.

Lemma 2.2. *For any symmetric scoring function $R : \mathcal{A}^* \times \mathcal{A}^* \rightarrow \mathbb{R}$, the following convergence bound applies,*

$$\left| \lambda_R - E \left[\frac{L_{n,R}}{n} \right] \right| \leq 3|R| \sqrt{\frac{\ln(ne)}{n}}.$$

Proof. To simplify the notation, let us write $\lambda_{n,R} = E[L_{n,R}]/n$. Let $m = kn$ for some $k \in \mathbb{N}$, and let $\mathcal{P}^{m,n}$ be the set of pairs of partitions of the integer interval $[1, m]$ into $2k$ pieces for which the sum of the lengths of the i -th pieces is always n . In other words,

$$\mathbf{p} = (i_0, i_1, \dots, i_{2k}, j_0, j_1, \dots, j_{2k})$$

is in $\mathcal{P}^{m,n}$ if

$$0 = i_0 < i_1 < \dots < i_{2k} = m,$$

$$0 = j_0 < j_1 < \dots < j_{2k} = m, \quad \text{and}$$

$$i_\ell - i_{\ell-1} + j_\ell - j_{\ell-1} = n, \quad \forall \ell \in [1, 2k].$$

For a partition $\mathbf{p} \in \mathcal{P}^{n,m}$, let $L_{m,R}^{\mathbf{p}}$ denote the optimal alignment score of $X_{[1,n]}$ and $Y_{[1,n]}$ relative to R under the extra constraint that the l -th pieces of the two partitions are aligned with each other, hence imposing that $X_{i_{l-1}+1} \dots X_{i_l}$ be aligned with $Y_{j_{l-1}+1} \dots Y_{j_l}$ for $l = 1, \dots, 2k$. In other words, we have

$$(2.3) \quad L_{m,R}^{\mathbf{p}} = \sum_{l=2}^{2k} R(X_{i_{l-1}+1} \dots X_{i_l}, Y_{j_{l-1}+1} \dots Y_{j_l}).$$

We can apply Azuma-Hoeffding to our constrained optimal alignment score $L_{m,R}^{\mathbf{p}}$ to justify that for any constant $\epsilon > 0$,

$$(2.4) \quad P(L_{m,R}^{\mathbf{p}} - E[L_{m,R}^{\mathbf{p}}] \geq \epsilon m) \leq \exp \left(-\frac{\epsilon^2 \cdot m}{2|R|^2} \right).$$

The optimal alignment score $L_{m,R}$ is not always equal to one of the the constrained alignment scores $L_{m,R}^{\mathbf{p}}$ however we can argue that it is not far from this. In fact, it is

not hard to see that for some partition \mathbf{p}

$$(2.5) \quad |L_{m,R} - L_{m,R}^{\mathbf{p}}| \leq 4k|R|.$$

Therefore, if the alignment score $L_{m,R}$ is to exceed a given benchmark, at least one of the constrained scores $L_{m,R}^{\mathbf{p}}$ must exceed this benchmark shifted by the correction term (2.5). This implies

$$(2.6) \quad \mathbb{P}[L_{m,R} \geq n\lambda_{n,R}k + \epsilon m] \leq \sum_{\mathbf{p} \in \mathcal{P}^{m,n}} \mathbb{P}[L_{m,R}^{\mathbf{p}} \geq n\lambda_{n,R}m + \epsilon m - 4k|R|].$$

We claim that by symmetry of R , we have

$$(2.7) \quad \mathbb{E}[L_{m,R}^{\mathbf{p}}] \leq n\lambda_{n,R}k, \quad \forall \mathbf{p} \in \mathcal{P}^{n,m}.$$

Our claim holds for two reasons: Firstly, $i_l - i_{l-1} + j_l - j_{l-1} = n$ implies $i_l < i_{l-1} + n$ and $j_l < j_{l-1} + n$ and

$$\begin{aligned} R(X_{i_{l-1}+1} \dots X_{i_l}, Y_{j_{l-1}+1} \dots Y_{j_l}) + R(X_{i_l+1} \dots X_{i_{l-1}+n}, Y_{j_l+1} \dots Y_{j_{l-1}+n}) \\ \leq R(X_{i_{l-1}+1} \dots X_{i_{l-1}+n}, Y_{j_{l-1}+1} \dots Y_{j_{l-1}+n}). \end{aligned}$$

Taking expectations on both sides, we find

$$(2.8) \quad \begin{aligned} \mathbb{E}[R(X_{i_{l-1}+1} \dots X_{i_l}, Y_{j_{l-1}+1} \dots Y_{j_l})] + \mathbb{E}[R(X_{i_l+1} \dots X_{i_{l-1}+n}, Y_{j_l+1} \dots Y_{j_{l-1}+n})] \\ \leq \mathbb{E}[R(X_{i_{l-1}+1} \dots X_{i_{l-1}+n}, Y_{j_{l-1}+1} \dots Y_{j_{l-1}+n})]. \end{aligned}$$

Secondly, the crucial assumption that R be symmetric implies that the two terms on the left-hand side of (2.8) are equal, thus yielding

$$(2.9) \quad \begin{aligned} 2\mathbb{E}[R(X_{i_{l-1}+1} \dots X_{i_l}, Y_{j_{l-1}+1} \dots Y_{j_l})] &\leq \mathbb{E}[R(X_{i_{l-1}+1} \dots X_{i_{l-1}+n}, Y_{j_{l-1}+1} \dots Y_{j_{l-1}+n})] \\ &= \mathbb{E}[R(X_1 \dots X_n, Y_1 \dots Y_n)] \\ &= n\lambda_{n,R}. \end{aligned}$$

Taking the expectation on both sides of (2.3) and applying (2.9) to each term on the right-hand side yields the claimed inequality, (2.7).

Substitution of (2.7) into (2.6) now yields

$$(2.10) \quad \mathbb{P}[L_{m,R} \geq n\lambda_{n,R}k + \epsilon m] \leq \sum_{\mathbf{p} \in \mathcal{P}^{m,n}} \mathbb{P}[L_{m,R}^{\mathbf{p}} \geq \mathbb{E}[L_{m,R}^{\mathbf{p}}] + \epsilon m - 4k|R|].$$

Using (2.4) and the fact that $\mathcal{P}^{n,m}$ has fewer than $\binom{m}{k}^2$ elements yields that for large n and k ,

$$(2.11) \quad \mathbb{P}[L_{m,R} \geq n\lambda_{n,R}k + \epsilon m] \leq \binom{m}{k}^2 \exp\left(-\frac{(\epsilon - 4|R|/n)^2 \cdot m}{2|R|^2}\right).$$

Let Z be a binomial variable with parameters m and $p = 1/n$, so that we have

$$\mathbb{P}[Z = k] = \binom{m}{k} \left(\frac{1}{n}\right)^k \cdot \left(\frac{n-1}{n}\right)^{m-k} \leq 1,$$

and hence,

$$(2.12) \quad \binom{m}{k} \leq n^k \cdot \left(\frac{1}{1 - \frac{1}{n}}\right)^{k(n-1)} \leq (e \cdot n)^k, \quad (n \gg 1).$$

Substituting (2.12) into (2.11), we find that for large n ,

$$(2.13) \quad \mathbb{P}\left[\frac{L_{m,R}}{m} \geq \lambda_{n,R} + \epsilon\right] \leq \exp\left(k \left[2 \ln(e \cdot n) - \frac{(\epsilon - 4|R|/n)^2 \cdot n}{4|R|^2}\right]\right).$$

The key now is to let k tend to infinity. In doing so, we know on the one hand that $L_{m,R}/m \rightarrow \lambda_R$, and on the other that the the right-hand side of (2.13) converges either to 0 or $+\infty$. It does converge to 0 only if

$$2 \ln(e \cdot n) - \frac{(\epsilon - 4|R|/n)^2 \cdot n}{4|R|^2} < 0$$

which is certainly satisfied if n is chosen large enough ($n > 10$ suffices) and

$$\epsilon = 3|R|\sqrt{\frac{\ln(ne)}{n}}.$$

Therefore, we find

$$\mathbb{P}\left[\lambda_R \geq \lambda_{n,R} + 3|R|\sqrt{\frac{\ln(ne)}{n}}\right] = 0,$$

and since λ_R is a constant, and similarly $\lambda_{n,R}$, we actually deduce that

$$\lambda_R \leq \lambda_{n,R} + 3|R|\sqrt{\frac{\ln(ne)}{n}}.$$

On the other hand, we also know from (2.1) that $\lambda_n/n \leq \lambda_R$, thus concluding the proof. \square

Lemma 2.3. *Let $R : \mathcal{A}^* \times \mathcal{A}^* \rightarrow \mathbb{R}$ be a symmetric scoring function and let $A^n(R)$ denote the event*

$$A^n(R) = \left\{ \left| \lambda_R - \frac{L_{n,R}}{n} \right| \leq 5|R| \sqrt{\frac{\ln(ne)}{n}} \right\}.$$

Then for large n ,

$$\mathbb{P}[A^n(R)] \geq 1 - n^{-\ln n}.$$

Proof. This follows by combining (2.2) with $\epsilon = 2|R|$, Lemma 2.2, Theorem 2.1 and Lemma 2.1. \square

The next result is about the convergence of convex compact sets.

Lemma 2.4. *Let $(A_n)_{n \in \mathbb{N}}$ be a sequence of random compact convex sets in \mathbb{R}^2 such that for any linear form $f \in (\mathbb{R}^2)^*$ there exists a deterministic constant $\xi_f \in \mathbb{R}$ for which*

$$\max_{(x,y) \in A_n} f(x,y) \xrightarrow{n \rightarrow \infty} \xi_f, \quad a.s.$$

Then there exists a deterministic compact convex set $A \subset \mathbb{R}^2$ for which

$$d_H(A_n, A) \xrightarrow{n \rightarrow \infty} 0, \quad a.s.,$$

where d_H is the Hausdorff distance.

Proof. Let F be a dense countable subset of the unit sphere in $(\mathbb{R}^2)^*$. Then

$$A := \{(x, y) : f(x, y) \leq \xi_f, \forall f \in (\mathbb{R}^2)^*\} = \{(x, y) : f(x, y) \leq \xi_f, \forall f \in F\}.$$

Furthermore, A is compact and convex, the condition of the lemma implies that

$$(2.14) \quad \mathbb{P} \left[\max_{(x,y) \in A_n} f(x,y) \xrightarrow{n \rightarrow \infty} \xi_f, \forall f \in F \right] = 1,$$

and we have

$$(2.15) \quad \max_{(x,y) \in A} f(x,y) = \xi_f, \quad \forall f \in F.$$

Suppose it is not the case that $d_H(A_n, A) \rightarrow 0$ almost surely. Then there exists $\delta > 0$ and a set $\mathcal{E} \subset \Omega$ such that $\mathbb{P}[\mathcal{E}] > 0$ and $\forall \omega \in \mathcal{E}$ there exists a sequence of points $(\alpha_n(\omega))_{n \in \mathbb{N}}$ such that $\alpha_n(\omega) \in A_n(\omega)$ and

$$d(\alpha_n(\omega), A) := \min_{\beta \in A} d(\alpha_n(\omega), \beta) \geq \delta.$$

Since all sets $A_n(\omega)$ are contained in some large closed box, there exists a convergent subsequence $(\alpha_{n_k}(\omega))_{k \in \mathbb{N}} \rightarrow \alpha(\omega)$. The continuity of the function $\alpha \mapsto d(\alpha, A)$ implies that we have $d(\alpha(\omega), A) \geq \delta > 0$, and in particular that $\alpha(\omega) \notin A$. By virtue of the Hahn-Banach separation theorem, there exists $g_\omega \in (\mathbb{R}^2)^*$ such that $A \subset \{(x, y) : g_\omega(x, y) \leq \max_{(s,t) \in A} g_\omega(s, t)\}$ and $g_\omega(\alpha) > \max_{(s,t) \in A} g_\omega(s, t) + \epsilon$ for some $\epsilon > 0$. Let $(f_\ell)_{\ell \in \mathbb{N}} \subset F$ be a sequence such that $f_\ell \rightarrow g_\omega$ in the weak topology. By (2.15), we have $A \subset \{(x, y) : f_\ell(x, y) \leq \xi_{f_\ell}\}$, and for ℓ large enough it is the case that $f_\ell(\alpha) > \xi_{f_\ell} + 2\epsilon/3$. If it were now the case that

$$(2.16) \quad \max_{(x,y) \in A_n(\omega)} f_\ell(x, y) \rightarrow \xi_{f_\ell},$$

then for large enough n ,

$$f_\ell(\alpha(\omega)) > \xi_{f_\ell} + 2\epsilon/3 > \max_{(x,y) \in A_{n_k}(\omega)} f_\ell(x, y) + \epsilon/3 \geq f_\ell(\alpha_{n_k}(\omega)) + \epsilon/3.$$

But this is a contradiction, since by continuity of f_ℓ , we have $f_\ell(\alpha_{n_k}(\omega)) \rightarrow f_\ell(\alpha(\omega))$. We conclude that for each $\omega \in \mathcal{E}$ there exists $f_\ell \in F$ for which (2.16) does not apply, and since $P[\mathcal{E}] > 0$, this contradicts (2.14). \square

Lemma 2.5. *Let S, T be two symmetric scoring functions on $\mathcal{A}^* \times \mathcal{A}^*$. If the optimization problem $\max_{(x,y) \in \text{SET}_{S,T}} x$ has a unique maximizer (x_0, y_0) , then*

$$(2.17) \quad \max_{\pi \in \nu_{n,S}^*} \frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} \xrightarrow{n \rightarrow \infty} y_0, \quad a.s.,$$

$$(2.18) \quad \min_{\pi \in \nu_{n,S}^*} \frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} \xrightarrow{n \rightarrow \infty} y_0, \quad a.s.$$

Proof. By virtue of (1.3) and Lemma 2.4, $d_H(\text{SET}_{S,T}^n, \text{SET}_{S,T}) \rightarrow 0$ almost surely. Keeping in mind (1.6) and (1.7), taking any convergent subsequence $((x_{n_\ell}, y_{n_\ell}))_{\ell \in \mathbb{N}}$ of a sequence $((x_n, y_n))_{n \in \mathbb{N}}$ of maximizers

$$(2.19) \quad (x_n, y_n) \in \arg \max \left\{ y : (x, y) \in \text{SET}_{S,T}^n, x = \frac{L_{n,S}}{n} \right\},$$

and writing $(x^*, y^*) = \lim_{\ell \rightarrow \infty} (x_{n_\ell}, y_{n_\ell})$, we have $x^* = x_0$ almost surely (by virtue of (1.3)), and $(x^*, y^*) \in \text{SET}_{S,T}$ almost surely. By the assumptions of the lemma, we thus have $(x^*, y^*) = (x_0, y_0)$. Furthermore, a convergent subsequence of $((x_n, y_n))_{n \in \mathbb{N}}$ always exists, since all sets $\text{SET}_{S,T}^n$ are contained in a compact box, and the argument above

shows that (x_0, y_0) is the only accumulation point. Therefore, $(x_n, y_n) \rightarrow (x_0, y_0)$ almost surely, and since the choice of (x_n, y_n) among the maximisers of (2.19) was arbitrary, (2.17) and (2.18) both follow. \square

Lemma 2.6. *Let $K \subset \mathbb{R}^2$ be a deterministic convex compact set. Then the maximizer*

$$(x_0, y_0) = \arg \max_{(x,y) \in K} ax + by$$

is unique for all but countable many points (a, b) on the unit sphere in \mathbb{R}^2 . Furthermore, if (a, b) is chosen uniformly at random from the unit sphere in \mathbb{R}^2 , then

$$(2.20) \quad \mathbb{P} [\kappa(\partial K, (x_0, y_0)) \leq k] \leq \frac{k \cdot l}{2\pi},$$

where $\kappa(\partial K, (x_0, y_0))$ is the curvature of the boundary of K at the point (x_0, y_0) , and where l denotes the length of the boundary of K .

Proof. The first part of the lemma is well known. The mapping

$$H : (a, b) \mapsto (x_0, y_0) := \arg \max_{(x,y) \in K} ax + by$$

is thus well defined for all but a countable number of points (a, b) on the unit circle. If the interior of K is empty, then K lives on a line segment. The maximiser (x_0, y_0) is then one of the two endpoints of this segment for almost all (a, b) , and since the curvature is infinite at these points, the claim of the lemma is trivially true.

If K has nonempty interior, then its boundary ∂K is locally the graph of a convex function, and hence it is continuous. Spherical projection with respect to an interior point defines a parametrization $u(\theta)$ of ∂K , with θ running on the unit circle. Since the boundary ∂K is locally the graph of a convex function, $u(\theta)$ is differentiable everywhere except at a countable number of points, and it is twice differentiable everywhere except on a set of Lebesgue measure 0, see e.g. [3, Theorem 1, page 242]. The length $l = \int_0^{2\pi} \|du(\theta)/d\theta\|_2 d\theta$ is thus well defined and finite, and so is the reparametrization $c(t)$ of $u(\theta)$ with respect to the length $t = \int_0^\theta \|du(\tau)/d\tau\|_2 d\tau$. Furthermore, we have $\|\dot{c}(t)\|_2 = 1$ for all $t \in [0, l]$, and $\ddot{c}(t)$ is defined except on a Lebesgue-null set. Let A be the subset of $t \in [0, l]$ where $\dot{c}(t)$ is defined, and B the subset where $\ddot{c}(t)$ is defined. Without loss of generality, we may assume that the orientation of the curve $c(t)$ is positive, so that $G(t) = i\dot{c}(t)$ is the unit normal vector to K at $c(t)$ (orthogonal to $\dot{c}(t)$ and

pointing away from K). This defines a mapping $t \mapsto G(t)$ from A to the unit circle. We make the following two observations:

- (a) $\kappa(t) := \kappa(\partial K, c(t)) = \|\ddot{c}(t)\|_2 = \|\dot{G}(t)\|_2$ equals the curvature of ∂K at $c(t)$.
- (b) Given (a, b) on the unit circle, if $c(t) = \arg \max_{(x,y) \in K} ax + by$ for some $t \in A$, and if this is the unique maximizer, then $(a, b) = G(t)$.

Let $T := \{t \in [0, l] : \kappa(t) \leq k\}$. The fact that $G(t)$ is defined at all points where $\kappa(t)$ is defined combined with Observations (a) and (b) imply that

$$\begin{aligned} \text{P}[\kappa(t) \leq k] &= \text{P}[G(T)] \\ (2.21) \quad &= \int_T |\dot{G}(t)| \frac{dt}{2\pi} \\ &= \int_T k(t) \frac{dt}{2\pi}, \end{aligned}$$

$$(2.22) \quad \leq \frac{k \cdot l}{2\pi}.$$

Equation (2.22) establishes the claim (2.20) of the lemma. The only nontrivial step that needs further explanation is (2.21). Let $g : [0, l] \rightarrow [0, 2\pi]$ be such that $G(t) = \exp(i g(t))$. Then $g(t)$ is well defined except at a countable number of points. By convexity of K , $g(t)$ is a non-decreasing function, and without loss of generality we may assume that it is right continuous. Equation (2.21) can thus be reformulated as follows,

$$(2.23) \quad \mu[g(T)] = \int_T \frac{\dot{g}(t)}{2\pi} dt,$$

where μ is the uniform probability measure on the interval $[0, 2\pi]$. If g is smooth and increasing, then (2.23) is simply a change of variable formula. In the general case we can approximate using smooth functions. Thus take a standard mollifier ϕ_ϵ and $g_{\epsilon,\delta} = (g + \delta h) \star \phi_\epsilon$ where $h(x) = x$. The rationale for taking $g + \delta h$ is to render the derivative positive and g_ϵ increasing. Equation (2.23) is true for $g_{\epsilon,\delta}$, and its general validity is obtained by first passing ϵ to zero, followed by δ . \square

3. PROOFS OF THE MAIN THEOREMS

3.1. Proof of Theorem 1.1.

Proof. Let $R : \mathcal{A}^* \times \mathcal{A}^* \rightarrow \mathbb{R}$ be a symmetric scoring function, and consider the event

$$A^n(R) = \left\{ \left| \lambda_R - \frac{L_{n,R}}{n} \right| \leq \frac{5|R|\sqrt{\ln(en)}}{\sqrt{n}} \right\}.$$

It follows from Lemma 2.2 and Theorem 2.1 that

$$\begin{aligned} \mathbb{P}[A^n(R)] &\geq \mathbb{P} \left[\left| \frac{L_{n,R}}{n} - \frac{E[L_{n,R}]}{n} \right| \leq 2|R|\sqrt{\ln(en)/n}, \left| \frac{E[L_{n,R}]}{n} - \lambda_R \right| \leq 3|R|\sqrt{\ln(en)/n} \right] \\ &\geq 1 - n^{-\ln n}, \quad \forall n \gg 1. \end{aligned}$$

By the assumptions of the theorem, $x_0 = \lambda_S$ and

$$(3.1) \quad \kappa(\partial \text{SET}_{S,T}, (x_0, y_0)) \geq k > 0.$$

For small $\epsilon > 0$ the point $P_\epsilon := (x_\epsilon, y_\epsilon)$ on the boundary $\partial \text{SET}_{S,T}$ with y -coordinate $y_\epsilon := y_0 + \epsilon/k$ nearest to (x_0, y_0) is well defined. Choose a_ϵ such that the linear form $f_{(1,a_\epsilon)} : (x, y) \mapsto x + a_\epsilon y$ has its maximizer over the set $\text{SET}_{S,T}$ at P_ϵ . This implies that for any $(x, y) \in \text{SET}_{S,T}$,

$$(3.2) \quad x + a_\epsilon y \leq x_\epsilon + a_\epsilon y_\epsilon.$$

The curvature condition (3.1) implies that for all ϵ small enough,

$$x_\epsilon \leq x_0 - \frac{\epsilon^2}{3}.$$

Combined with (3.2) for $(x, y) = (x_0, y_0)$, this yields $(x - x_\epsilon) + a_\epsilon(y - y_\epsilon) \leq 0$, and since furthermore $x_0 > x_\epsilon$, it follows that

$$(3.3) \quad a_\epsilon \geq \frac{x_0 - x_\epsilon}{y_\epsilon - y_0} \geq \frac{\epsilon k}{3}.$$

If $A^n(S)$ holds, then for any optimal alignment π relative to S we have

$$(3.4) \quad \left| \frac{S_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - x_0 \right| \leq \frac{5|S|\sqrt{\ln(en)}}{\sqrt{n}},$$

and similarly, if the event $A^n(S + a_\epsilon T)$ holds, then

$$(3.5) \quad \left| \frac{L_{n,S+a_\epsilon T}}{n} - \lambda_{S+a_\epsilon T} \right| \leq \frac{5|S|\sqrt{\ln(en)}}{\sqrt{n}}.$$

On the other hand,

$$\lambda_{S+a_\epsilon T} = \max_{(x,y) \in \text{SET}_{S,T}} f_{(1,a_\epsilon)}(x,y) = x_\epsilon + a_\epsilon y_\epsilon,$$

and substituted into (3.5) this yields

$$(3.6) \quad \left| \frac{L_{n,S+a_\epsilon T}}{n} - (x_\epsilon + a_\epsilon y_\epsilon) \right| \leq \frac{5|S + a_\epsilon T|\sqrt{\ln(en)}}{\sqrt{n}}.$$

Next, for any optimal alignment π relative to S , we have

$$\frac{(S + a_\epsilon T)_\pi(X_{[1,n]}, Y_{[1,n]})}{n} \leq \frac{L_{n,S+a_\epsilon T}}{n} \stackrel{(3.6)}{\leq} x_\epsilon + a_\epsilon y_\epsilon + \frac{5|S + a_\epsilon T|\sqrt{\ln(en)}}{\sqrt{n}}.$$

It now follows from (3.4) that

$$a_\epsilon \left(\frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - y_0 \right) \leq x_\epsilon - x_0 + a_\epsilon(y_\epsilon - y_0) + \frac{5(|S + a_\epsilon T| + |S|)\sqrt{\ln(en)}}{\sqrt{n}},$$

and since $x_\epsilon - x_0 \leq 0 < a_\epsilon$, this finally yields that for large n and small $\epsilon > 0$,

$$\frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - y_0 \leq \frac{\epsilon}{k} + \frac{5(|S + a_\epsilon T| + |S|)\sqrt{\ln(en)}}{a_\epsilon \sqrt{n}} \leq \frac{5(2|S| + a_\epsilon |T|)\sqrt{\ln(en)}}{a_\epsilon \sqrt{n}}.$$

In combination with (3.3) this yields

$$\frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - y_0 \leq \frac{\epsilon}{k} + \frac{5(6|S| + \epsilon k |T|)}{\epsilon k} \sqrt{\frac{\ln(en)}{n}}.$$

For large n , we can minimize the right-hand side over ϵ , yielding

$$\frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - y_0 \leq \frac{5|T| + 2\sqrt{30|S|}}{k} \left(\frac{\ln(en)}{n} \right)^{1/4}$$

By changing the scoring function T to $-T$, an analogous argument also shows that

$$- \left(\frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - y_0 \right) \leq \frac{5|T| + 2\sqrt{30|S|}}{k} \left(\frac{\ln(en)}{n} \right)^{1/4},$$

and hence,

$$(3.7) \quad \left| \frac{T_\pi(X_{[1,n]}, Y_{[1,n]})}{n} - y_0 \right| \leq \frac{5|T| + 2\sqrt{30|S|}}{k} \left(\frac{\ln(en)}{n} \right)^{1/4}.$$

We conclude that if all of the events $A^n(S)$ and $A^n(S + a_\epsilon T)$ and $A^n(S - a_\epsilon T)$ hold, then (3.7) applies, and since the probability that any individual event fails to hold is bounded by $n^{-\ln n}$, the claim of the theorem follows. \square

3.2. Proof of Theorem 1.2.

Proof. Let $V = \arccos\langle S, T \rangle_F$, where $\langle \cdot, \cdot \rangle_F$ is the inner product on the space of symmetric scoring functions that corresponds to the Frobenius norm. Then V is uniformly distributed on $[-\pi/2, \pi/2]$. Let T_1 be the Gram-Schmidt orthogonalization of T with respect to $S_1 := S$, and let U be a uniform random variable on $[0, 2\pi]$, independent of S and T , and hence also of V , and let us define $(S_2, T_2) = \Phi(S_1, T_1)$, where Φ is the rotation

$$\begin{aligned} \Phi : \mathbb{R}^2 &\rightarrow \mathbb{R}^2, \\ (x, y) &\mapsto (\cos(U)x + \sin(U)y, -\sin(U)x + \cos(U)y) \end{aligned}$$

by the angle U . It is easy to see that $\text{SET}_{S_2, T_2} = \Phi(\text{SET}_{S_1, T_1})$, and that under Φ^{-1} , the point where SET_{S_2, T_2} has a point of maximal first coordinate corresponds to the point where the random linear form $f : (x, y) \mapsto \cos(U)x + \sin(U)y$ takes a maximum value on SET_{S_1, T_1} . Furthermore, since Φ is angle-preserving, the curvature κ_1 of $\partial\text{SET}_{S_2, T_2}$ and $\partial\text{SET}_{S_1, T_1}$ at these points is also the same. Lemma 2.6 applies, and we have $\mathbb{P}[\kappa_1 \leq k] \leq k \cdot l / (2\pi)$, where l is the length of the boundary of SET_{S_1, T_1} . Since the scoring functions under considerations have unit norm, the rescaled alignment score cannot exceed 2, implying that $l \leq 8$ and

$$(3.8) \quad \mathbb{P}[\kappa_1 \leq k] \leq \frac{4k}{\pi}.$$

It remains to relate κ_1 to the curvature κ of $\partial\text{SET}_{S, T}$ at the point where its first coordinate is maximized. Since $\text{SET}_{S, T} = \Psi(\text{SET}_{S_1, T_1})$, where Ψ is the linear transformation

$$\begin{aligned} \Psi : \mathbb{R}^2 &\rightarrow \mathbb{R}^2, \\ (x, y) &\mapsto (x, \cos(V)x + \sin(V)y), \end{aligned}$$

we have $\kappa = \kappa_1 / |\sin V| \geq \kappa_1$, so that

$$P[\kappa < k] \leq P[\kappa_2 < k] \leq \frac{4k}{\pi},$$

as claimed in the statement of the theorem. \square

REFERENCES

- [1] K. Azuma. Weighted sums of certain dependent random variables. *Tohoku Math. J.*, 19:357–367, 1967.
- [2] Václav Chvatal and David Sankoff. Longest common subsequences of two random sequences. *J. Appl. Probability*, 12:306–315, 1975.
- [3] Lawrence C. Evans and Ronald F. Gariepy. *Measure theory and fine properties of functions*. Studies in Advanced Mathematics. CRC Press, Boca Raton, FL, 1992.
- [4] Raphael Hauser and Heinrich Matzinger. Distribution of aligned letter pairs in optimal alignments of random sequences. *arXiv:1211.5491*, 2013.
- [5] S.B. Needleman and C.D. Wunsch. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology*, 48:443–453, 1970.
- [6] M.S. Waterman and M. Vingron. Sequence comparison significance and poisson approximation. *Statistical Science*, 9(3):367–381, 1994.

RAPHAEL HAUSER, MATHEMATICAL INSTITUTE, UNIVERSITY OF OXFORD, RADCLIFFE OBSERVATORY QUARTER, WOODSTOCK ROAD, OXFORD OX2 6GG, UNITED KINGDOM, AND PEMBROKE COLLEGE, ST ALDATES, OXFORD, OX1 1DW, UNITED KINGDOM.

E-mail address: hauser@maths.ox.ac.uk

HEINRICH MATZINGER, SCHOOL OF MATHEMATICS, GEORGIA INSTITUTE OF TECHNOLOGY, 686 CHERRY STREET, ATLANTA, GA 30332-0160 USA. CORRESPONDING AUTHOR.

E-mail address: matzi@math.gatech.edu

IONEL POPESCU, SCHOOL OF MATHEMATICS, GEORGIA INSTITUTE OF TECHNOLOGY, 686 CHERRY STREET, ATLANTA, GA 30332, USA, AND “SIMION STOILOW” INSTITUTE OF MATHEMATICS OF ROMANIAN ACADEMY, 21 CALEA GRIVIȚEI, BUCHAREST, ROMANIA.

E-mail address: ipopescu@math.gatech.edu, ionel.popescu@imar.ro